

An Efficient Privacy Attack on MyAnimeList's Affinity Oracle: Technical Notes on Score Estimation

Stephen Huan

March 7, 2021

For context (and code) see this [GitHub repository](#).

1 Pearson's Correlation

Definition 1.1. *Pearson's correlation* between vectors \mathbf{u} and \mathbf{v} is defined as follows: first, define $\mathbf{u}' = \mathbf{u} - \bar{\mathbf{u}}$ where $\bar{\mathbf{u}}$ is the constant vector equal to the sample mean of \mathbf{u} , and define \mathbf{v}' similarly. Then Pearson's correlation is simply the cosine similarity:

$$\rho(\mathbf{u}, \mathbf{v}) = \cos \alpha = \frac{\mathbf{u}' \cdot \mathbf{v}'}{|\mathbf{u}'||\mathbf{v}'|}$$

Theorem 1.1. *Pearson's correlation is invariant to a transformation of the form $a\mathbf{u} + \mathbf{b}$, for any positive real number a and constant vector \mathbf{b} .*

Proof. We start with the constant vector, which is destroyed when taking the mean:

$$\begin{aligned} \mathbf{u}' &= \mathbf{u} - \mathbb{E}[\mathbf{u}] && \text{definition} \\ (a\mathbf{u} + \mathbf{b})' &= a\mathbf{u} + \mathbf{b} - (a\mathbb{E}[\mathbf{u}] + \mathbf{b}) && \text{linearity of expectation} \\ &= a(\mathbf{u} - \mathbb{E}[\mathbf{u}]) \\ &= a\mathbf{u}' \end{aligned}$$

We now see what happens to this scaling by a after Pearson's correlation:

$$\rho(a\mathbf{u} + \mathbf{b}, \mathbf{v}) = \frac{(a\mathbf{u}') \cdot \mathbf{v}'}{|a\mathbf{u}'||\mathbf{v}'|} = \frac{a(\mathbf{u}' \cdot \mathbf{v}')}{a|\mathbf{u}'||\mathbf{v}'|} = \rho(\mathbf{u}, \mathbf{v})$$

If $a = 0$ then dividing by the magnitude of the zero vector is ill-defined, and if $a < 0$ then $|a\mathbf{u}'|$ becomes $|a||\mathbf{u}'|$, so the sign of the correlation switches because of the a on the numerator. Thus, a must be positive for the correlation to be preserved. \square

Corollary 1.1.1. *We can't tell the difference between \mathbf{u} and $a\mathbf{u} + \mathbf{b}$ by looking at Pearson's correlation without using additional information.*

2 The Attack

Suppose \mathbf{u} is an unknown vector and our only method of obtaining information about \mathbf{u} is through querying some *affinity oracle* which tells us $\rho(\mathbf{u}, \mathbf{v})$ for some \mathbf{v} that we pick. We now show how to pick \mathbf{v} in such a way that \mathbf{u} can be determined, up to a positive scaling and constant translation as mentioned previously.

The basic idea is going to be to pick \mathbf{v} such that \mathbf{v}' is as close to a basis vector as possible, e.g. something like $[1 \ 0 \ 0 \ \dots \ 0]^T$. In that case when we compute $\rho(\mathbf{u}, \mathbf{v})$, it equals $\mathbf{u}' \cdot \mathbf{v}'$ divided by the magnitude of \mathbf{u} , which we can ignore since we can't tell the difference between \mathbf{u} and $a\mathbf{u}$ anyways. Since \mathbf{v}' is all zeros except for one 1, $\mathbf{u}' \cdot \mathbf{v}' = u'_1$, telling us the first element of \mathbf{u}' directly. The first snag is that it is not possible for such a \mathbf{v}' to exist, since the sum of any \mathbf{x}' must be 0 and \mathbf{v}' sums to 1.

Theorem 2.1. *The sum of the elements in \mathbf{x}' is 0 for any \mathbf{x} .*

Proof.

$$\mathbb{E}[\mathbf{x}'] = \mathbb{E}[\mathbf{x} - \mathbb{E}[\mathbf{x}]] = \mathbb{E}[\mathbf{x}] - \mathbb{E}[\mathbf{x}] = 0$$

□

Thus, we need a different pick for \mathbf{v}' . One simple fix is to introduce a -1 to cancel the existing 1, i.e. $[1 \ -1 \ 0 \ \dots \ 0]^T$. When we compute $\mathbf{u}' \cdot \mathbf{v}'$, we get that $\rho(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}' \cdot \mathbf{v}'}{|\mathbf{u}'||\mathbf{v}'|} = u'_1 - u'_2$ (where we arbitrarily let $|\mathbf{u}'| = \frac{1}{|\mathbf{v}'|}$ for simplicity). This is a linear equation in \mathbf{u} , so to solve for \mathbf{u} we just need a system. The first step is finding a \mathbf{v} that generates the \mathbf{v}' used above. That can be done by seeing a constant vector will map to the zero vector, so we can start with a vector of all 5's, and then put a 6 where the 1 in \mathbf{v}' should be and a 4 where the -1 should be. The average of \mathbf{v} will be 5, so \mathbf{v}' will map to zeros except the 6 goes to 1 and the 4 goes to -1, which is what we wanted. The second step is that we need a system of equations, so we need multiple (linearly-independent) \mathbf{v}' 's. The obvious choice is to "slide" the 1 over, so our next query is $[0 \ 1 \ -1 \ \dots \ 0]^T$, the next $[0 \ 0 \ 1 \ -1 \ \dots \ 0]^T$ and so on. Notice if \mathbf{u} has M elements, then we only have $M - 1$ versions of \mathbf{v} since we can't put a 1 at the last position (we need a -1 after it). This is a natural consequence of the fact that we are missing *two* variables: we can't tell the difference between \mathbf{u} and $a\mathbf{u} + \mathbf{b}$, so a and \mathbf{b} are free variables. By setting $|\mathbf{u}'|$ we lose a free variable, so we have one free variable left. Thus, our matrix must also have one free variable. We therefore solve the system by putting each variable in terms of the last variable. We have a system of the form:

$$\begin{bmatrix} 1 & -1 & & \\ & 1 & -1 & \\ & & 1 & -1 \end{bmatrix} \mathbf{u}' = \boldsymbol{\rho}$$

From the last row, $u'_{n-1} - u'_n = \rho_{n-1}$ so $u'_{n-1} = \rho_{n-1} + u'_n$. We can then add the last row to the second to last row, which cancels the -1 with the 1, yielding $[\dots \ 1 \ 0 \ -1]$. Thus, $u'_{n-2} = \rho_{n-2} + \rho_{n-1} + u'_n$. We can then add this row to row above it, and so on.

Thus, the value of \mathbf{u}_i is the sum $\rho_i + \rho_{i+1} + \dots + \rho_{n-1}$ which is a suffix sum that can be computed in $O(n)$ over each element in \mathbf{u}' . Now that we have \mathbf{u}' , we need to find a possible \mathbf{u} that transforms into it. This could be done directly, see [subsection 3.1](#) for a brute-force analysis, but we'll use a simple observation derived from [Theorem 2.1](#).

Theorem 2.2. $(\mathbf{x}')' = \mathbf{x}'$, i.e. \mathbf{x}' is a fixed point of $f(\mathbf{x}) = \mathbf{x} - \mathbb{E}[\mathbf{x}]$.

Proof. $(\mathbf{x}')' = \mathbf{x}' - \mathbb{E}[\mathbf{x}']$ by definition. By [Theorem 2.1](#), $\mathbb{E}[\mathbf{x}']$ is 0. Thus, $(\mathbf{x}')' = \mathbf{x}'$. \square

Recall we wanted to find a \mathbf{x} such that $\mathbf{x}' = \mathbf{u}'$. Since $(\mathbf{u}')' = \mathbf{u}'$, \mathbf{u}' is its own \mathbf{x} .

2.1 Determining an Valid Integer Form

Now that we have a possible \mathbf{u} , we need to determine the particular value of a and \mathbf{b} . One constraint is that each value in a score list must be an integer between 1 and 10, so $a\mathbf{u} + \mathbf{b} \in \mathbb{Z}^M$ and $1 \leq a\mathbf{u} + \mathbf{b} \leq 10$. Suppose we had a particular \mathbf{u} that satisfies those requirements. Then we could “normalize” \mathbf{u} by dividing by the greatest common divisor of the score differences, e.g. if $\mathbf{u} = (1 \ 3 \ 1 \ 3 \ 9 \ 9)$, we first sort the scores and remove duplicates into $\mathbf{u} = (1 \ 3 \ 9)$. We then notice $3 - 1 = 2$ and $9 - 3 = 6$, and $\text{gcd}(2, 6) = 2$. We can therefore transform \mathbf{u} by $\frac{1}{2}\mathbf{u} + \frac{1}{2}$ such that 1 goes to 1, 3 goes to 2, and 9 goes to 5. It is now not possible to multiply \mathbf{u} by a fraction and then shift \mathbf{u} to make all the elements integer because the elements will have different denominators. We sort and only consider the gcd of adjacent elements because if some factor d divides each index i to j , then it must also divide the difference of the element at i from the element at j since the difference is the sum of adjacent differences: $\mathbf{v}_j - \mathbf{v}_i = (\mathbf{v}_{i+1} - \mathbf{v}_i) + (\mathbf{v}_{i+2} - \mathbf{v}_{i+1}) + \dots + (\mathbf{v}_j - \mathbf{v}_{j-1})$. Thus, a is now constrained to \mathbb{Z} if \mathbf{u} is in this normalized form, and that constrains $\mathbf{b} \in \mathbb{Z}^M$ as well. There are now a finite number of possibilities for both, since $a > 0$ and the end result $a\mathbf{u} + \mathbf{b}$ must be between 1 and 10. Thus, if we have a *single* valid integer representation of \mathbf{u} , we can normalize it and then generate *all* valid possibilities.

2.2 Finding a Single Valid Integer Form

Recall that we have some $\mathbf{u} \in \mathbb{R}^M$ derived from solving a linear system. Since we obtain \mathbf{u} by computing the suffix sum of the vector of correlations ρ , it will not be integer. We therefore need to find some \mathbf{v} such that $\mathbf{v} \in \mathbb{Z}^M$, $1 \leq \mathbf{v} \leq 10$, and $\mathbf{v} = a\mathbf{u} + \mathbf{b}$ for some positive real number a and constant vector \mathbf{b} . Framed like this it seems like a search problem, and our search space is exponential: each score has 10 possibilities so a list of length M has 10^M possibilities. In order to cut down on the possibilities, we need to define a *canonical form*. \mathbf{u} and \mathbf{v} are *equivalent* if there exists a , \mathbf{b} such that $\mathbf{v} = a\mathbf{u} + \mathbf{b}$ (the interested reader can verify this defines an equivalence relation, see [subsection 3.2](#)). We now need a single, unique, vector to represent the equivalence class defined by all vectors equivalent to some vector \mathbf{u} . This is the *canonical form* for \mathbf{u} , or more precisely the canonical form for the equivalence class \mathbf{u} belongs to. We can do this simply by “removing” the parameters. We know what the maximum element in \mathbf{u} is. If we apply a transformation $a\mathbf{u} + \mathbf{b}$, then it will still be the maximum element by linearity.

Thus, the maximum and minimum elements give us hints about a . We can divide by $\max(\mathbf{u}) - \min(\mathbf{u})$ to force the range to be 1, and then shift the resulting vector such that the minimum element is 0 so that \mathbf{b} is fixed.

Theorem 2.3. *The canonical form for a vector \mathbf{u} is $\mathcal{C}(\mathbf{u}) = \frac{\mathbf{u} - \min(\mathbf{u})}{\max(\mathbf{u}) - \min(\mathbf{u})}$.*

Proof. We want to show that vectors \mathbf{u} and \mathbf{v} are equivalent if and only if their canonical forms are equal. We first prove the forward direction, if two vectors are equivalent, then their canonical forms are equal. By hypothesis, \mathbf{u} is equivalent to \mathbf{v} , so $\mathbf{v} = a\mathbf{u} + \mathbf{b}$.

$$\begin{aligned} \mathcal{C}(\mathbf{v}) &= \frac{\mathbf{v} - \min(\mathbf{v})}{\max(\mathbf{v}) - \min(\mathbf{v})} \\ &= \frac{a\mathbf{u} + \mathbf{b} - (a\min(\mathbf{u}) + \mathbf{b})}{(a\max(\mathbf{u}) + \mathbf{b}) - (a\min(\mathbf{u}) + \mathbf{b})} \\ &= \frac{a(\mathbf{u} - \min(\mathbf{u}))}{a(\max(\mathbf{u}) - \min(\mathbf{u}))} \\ &= \mathcal{C}(\mathbf{u}) \end{aligned}$$

We now prove the reverse direction, if $\mathcal{C}(\mathbf{u}) = \mathcal{C}(\mathbf{v})$ then \mathbf{u} and \mathbf{v} are equivalent. $\mathcal{C}(\mathbf{u})$ is of the form $a\mathbf{u} + \mathbf{b}$ where $a > 0$, so $\mathcal{C}(\mathbf{u}) = \mathcal{C}(\mathbf{v})$ implies $a_1\mathbf{u} + \mathbf{b}_1 = a_2\mathbf{v} + \mathbf{b}_2$. Thus, $\mathbf{v} = \frac{a_1}{a_2}\mathbf{u} + \frac{1}{a_2}[\mathbf{b}_1 - \mathbf{b}_2]$, and because $a_2 \neq 0$, $\frac{a_1}{a_2} > 0$, \mathbf{u} and \mathbf{v} are equivalent. \square

We can therefore transform the question “for a given vector \mathbf{u} , find a vector \mathbf{v} such that $\mathbf{v} \in \mathbb{Z}^M$, $1 \leq \mathbf{v} \leq 10$, and $\mathbf{v} = a\mathbf{u} + \mathbf{b}$ ” into “for a given vector \mathbf{u} , find a vector $\mathbf{v} \mid \mathbf{v} \in \mathbb{Z}^M, 1 \leq \mathbf{v} \leq 10$, and $\mathcal{C}(\mathbf{u}) = \mathcal{C}(\mathbf{v})$ ”. Now we don’t have to deal with \mathbf{u} at all, just its canonical form, which is helpful since we can construct \mathbf{v} from $\mathcal{C}(\mathbf{u})$. We first arbitrarily assume $\min(\mathbf{v}) = 0$; we can translate \mathbf{v} later. Given this assumption, $\mathcal{C}(\mathbf{v}) = \frac{\mathbf{v}}{\max(\mathbf{v})}$. Setting equal to $\mathcal{C}(\mathbf{u})$, $\frac{\mathbf{v}}{\max(\mathbf{v})} = \mathcal{C}(\mathbf{u})$ so $\mathbf{v} = \text{round}(\max(\mathbf{v})\mathcal{C}(\mathbf{u}))$ since we need \mathbf{v} to be integer. We can then iterate through the possible values for $\max(\mathbf{v})$, between 1 and 9 since we assume \mathbf{v} is not a constant vector and it is at most 9 greater than the minimum ($10 - 1 = 9$). We choose the maximum value whose computed \mathbf{v} is actually equivalent to \mathbf{u} , measured by the distance between the canonical forms, $|\mathcal{C}(\mathbf{u}) - \mathcal{C}(\mathbf{v})|$. Ideally this difference should be 0, but because our ρ is noisy (it is only given to 1 decimal place) we need to have some error tolerance. We are guaranteed that at least one $\mathcal{C}(\mathbf{v})$ is close to $\mathcal{C}(\mathbf{u})$ since there is a ground truth — we got \mathbf{u} through Pearson’s correlation against a ground truth score vector between 1 and 10.

2.3 Enumerating and Filtering Possibilities

Once we have this \mathbf{v} , we could still be off from the ground truth by a translation and a scaling. For example, if our vector has scores between 3–10, then we could subtract 2 from each element to yield an equivalent valid integer vector. With just Pearson’s, the integer constraint, and the 1–10 constraint we can’t distinguish between these two vectors. As mentioned in [subsection 2.1](#), if our vector is sufficiently normalized then a and \mathbf{b} are constrained to take integer values (and there is a finite number of pairs). We can

in fact implicitly normalize if in the iteration over $\max(\mathbf{v})$ we break ties by prioritizing the smaller value. Given this normalized \mathbf{v} , we now generate all valid a, \mathbf{b} pairs. We simply go through the possible values of a , between 1 and $\lfloor \frac{10}{\max(\mathbf{v})} \rfloor$ inclusive, and given a value of a the possible values of \mathbf{b} are those which map the smallest to least 1 and the maximum to at most 10, i.e. b can be between $1 - a \min(\mathbf{v})$ and $10 - a \max(\mathbf{v})$ inclusive.

Finally, we need to guess a particular possibility out of the possible pairs. MAL gives us the private lists's mean to two decimal places, so if the ground truth mean is μ and our mean is $\bar{\mathbf{v}}$, then we should pick a, \mathbf{b} such that our mean is as close to μ as possible, or when $|a\bar{\mathbf{v}} + |\mathbf{b}| - \mu|$ is minimized (this follows again from the linearity of expectation).

Another approach is to use maximum likelihood estimation, but that discussion is relegated to the appendix, [subsection 3.3](#) since knowing the sample mean μ is much better than knowing the distribution in practice.

3 Appendix

3.1 Undoing Mean Subtraction

We have a given vector \mathbf{u}' and want to find a vector \mathbf{u} such that $\mathbf{u} - \mathbb{E}[\mathbf{u}] = \mathbf{u}'$. We can solve for \mathbf{u} directly by setting up a linear system, although our linear system will necessarily have one free variable since the transformation is not invertible. The transformation $\mathbf{u} - \mathbb{E}[\mathbf{u}]$ will be shown to be represented as a matrix A times \mathbf{u} . This matrix is not invertible because the transformation is not surjective: our matrix is a mapping $\mathbb{R}^M \rightarrow \mathbb{R}^M$, and we cannot get vectors whose sum is not 0 by [Theorem 2.1](#). Thus, our function is not surjective and because the domain and co-domain are the same size, our function is not injective: \mathbf{x} and $\mathbf{x} + \mathbf{b}$ for some constant vector \mathbf{b} map to the same vector by [Theorem 1.1](#). Thus, our function cannot be invertible on some subspace of \mathbb{R}^M . This is all a long-winded way of saying we will have a free variable because we can't distinguish between \mathbf{x} and $\mathbf{x} + \mathbf{b}$. To write $f(\mathbf{u}) = \mathbf{u} - \mathbb{E}[\mathbf{u}]$ as a matrix, we simply write out f .

Taking the first element,

$$\begin{aligned} \mathbf{u}_1 - \frac{\mathbf{u}_1 + \mathbf{u}_2 + \dots + \mathbf{u}_M}{M} &= \mathbf{u}'_1 \\ \frac{(M-1)\mathbf{u}_1 - \mathbf{u}_2 - \dots - \mathbf{u}_M}{M} &= \mathbf{u}'_1 \\ (M-1)\mathbf{u}_1 - \mathbf{u}_2 - \dots - \mathbf{u}_M &= M\mathbf{u}'_1 \end{aligned}$$

Thus, we get a matrix of the form:

$$\begin{bmatrix} M-1 & -1 & -1 & -1 \\ -1 & M-1 & -1 & -1 \\ -1 & -1 & M-1 & -1 \\ -1 & -1 & -1 & M-1 \end{bmatrix}$$

If we subtract the second to last row from the last row, we get a row like $[0 \ 0 \ \dots \ -M \ M]$. We can then subtract the third to last row from the second to last row, getting a row like $[0 \ \dots \ -M \ M \ 0]$. Repeating the process, we get a matrix of the form:

$$\begin{bmatrix} -M & M & & \\ & -M & M & \\ & & -M & M \end{bmatrix}$$

where we discard the top row since we have a free variable. Amazingly, this matrix is a multiple of the 1 -1 matrix used in [section 2](#), so we have an equation of the form $-MA\mathbf{u} = M\mathcal{T}(\mathbf{u}')$, so $A\mathbf{u} = -\mathcal{T}(\mathbf{u}')$. $\mathcal{T}(\mathbf{u}')$ is from the fact that we subtracted adjacent rows to force the original matrix into a multiple of A , so $-\mathcal{T}(\mathbf{u}')$ looks like $[\mathbf{u}'_1 - \mathbf{u}'_2 \ \mathbf{u}'_2 - \mathbf{u}'_3 \ \dots \ \mathbf{u}'_{M-1} - \mathbf{u}'_M]^T$. When we take the suffix sum, we will get $\mathbf{u}'_i - \mathbf{u}'_M$, so our solution is $\mathbf{u}' - \mathbf{b}$. This could have been predicted from the start if we realized \mathbf{u}' was a constant difference from \mathbf{u} (they differ by $\mathbb{E}[\mathbf{u}]$) so \mathbf{u}' must have been a valid solution to the system to begin with. See [archive.py](#) for an implementation, although clearly we just went in a large circle (it's a very elegant circle, however).

3.2 Equivalence Relation Under Pearson's Correlation

Vectors \mathbf{u} and \mathbf{v} are *equivalent* if they are indistinguishable under Pearson's correlation, i.e. $\rho(\mathbf{u}, \mathbf{x}) = \rho(\mathbf{v}, \mathbf{x})$ for all \mathbf{x} . This is equivalent to saying there exists $a > 0$, \mathbf{b} such that $\mathbf{v} = a\mathbf{u} + \mathbf{b}$. Does this define an equivalence relation?

- Reflexivity: Pick $a = 1, \mathbf{b} = \mathbf{0}$ so \mathbf{v} is equivalent to itself.
- Symmetry. Suppose $\mathbf{v} = a\mathbf{u} + \mathbf{b}$. Then $\mathbf{u} = \frac{\mathbf{v}-\mathbf{b}}{a} = \frac{1}{a}\mathbf{v} - \frac{\mathbf{b}}{a}$, which is valid for $a \neq 0$.
- Transitivity. Suppose $\mathbf{v} = a\mathbf{u} + \mathbf{b}$, and $\mathbf{w} = c\mathbf{v} + \mathbf{d}$. Then $\mathbf{w} = c(a\mathbf{u} + \mathbf{b}) + \mathbf{d} = ca\mathbf{u} + [c\mathbf{b} + \mathbf{d}]$, which is valid since $ca > 0$.

3.3 Maximum Likelihood Estimation

Suppose we have two vectors, \mathbf{u} and \mathbf{v} . Which do we pick? We pick the one that is the most *likely* to occur, assuming we know the distribution. If we have a probability mass function f , $\mathbf{u} = (1 \ 2 \ 2 \ 4)$, and $\mathbf{v} = (1 \ 3 \ 5 \ 5)$, then the *likelihood* of \mathbf{u} occurring is $\mathcal{L}(\mathbf{u}) = f(1) \cdot f(2)^2 \cdot f(4)$ assuming each value of \mathbf{u} is identically and independently distributed (i.i.d.). Similarly, the likelihood of \mathbf{v} is $\mathcal{L}(\mathbf{v}) = f(1) \cdot f(3) \cdot f(5)^2$. We would choose \mathbf{u} or \mathbf{v} depending on which $\mathcal{L}(\mathbf{u})$ or $\mathcal{L}(\mathbf{v})$ was bigger. In general,

$$\mathcal{L}(\mathbf{x}) = \prod_{i=1}^M p(\mathbf{x}_i)$$

If we have the frequency of each score x , we can compute this more efficiently as:

$$\mathcal{L}(\mathbf{x}) = \prod_{x \in \mathbf{x}} p(x)^{\text{count}[x]}$$

Because \ln is a monotonic function, we can maximize $\ln \mathcal{L}(\mathbf{x})$ instead of $\mathcal{L}(\mathbf{x})$, which allows us to turn the product into a sum. This is useful because $\mathcal{L}(\mathbf{x})$ decreases exponentially, so we will run into precision issues without the log.

$$\ln \mathcal{L}(\mathbf{x}) = \sum_{x \in \mathbf{x}} \text{count}[x] \ln p(x)$$

Finally, to pick the a and \mathbf{b} which maximize the likelihood, we simply plug $a\mathbf{x} + \mathbf{b}$ in:

$$\ln \mathcal{L}(a\mathbf{x} + \mathbf{b}) = \sum_{x \in \mathbf{x}} \text{count}[x] \ln p(ax + b)$$